



**HAL**  
open science

# Simplitigs as an efficient and scalable representation of de Bruijn graphs

Karel Brinda, Michael Baym, Gregory Kucherov

► **To cite this version:**

Karel Brinda, Michael Baym, Gregory Kucherov. Simplitigs as an efficient and scalable representation of de Bruijn graphs. *Genome Biology*, 2021, 22 (96), 10.1101/2020.01.12.903443 . hal-03049400

**HAL Id: hal-03049400**

**<https://hal.science/hal-03049400>**

Submitted on 9 Dec 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# 1 Simplitigs as an efficient and scalable representation 2 of de Bruijn graphs

3 Karel Břinda<sup>1,2,\*</sup>, Michael Baym<sup>1,§</sup>, and Gregory Kucherov<sup>3,4,§</sup>

4 1 Department of Biomedical Informatics, Harvard Medical School, Boston, USA

5 2 Center for Communicable Disease Dynamics, Department of Epidemiology, Harvard T.H. Chan School of Public  
6 Health, Boston, USA

7 3 CNRS/LIGM Univ Gustave Eiffel, Marne-la-Vallée, France

8 4 Skolkovo Institute of Science and Technology, Moscow, Russia

9 \* The corresponding author: [karel.brinda@hms.harvard.edu](mailto:karel.brinda@hms.harvard.edu)

10 § These authors contributed equally to this work

## 11 Abstract

12 De Bruijn graphs play an essential role in computational biology. However, despite their widespread use, they lack a  
13 universal scalable representation suitable for different types of genomic data sets. Here, we introduce simplitigs as a  
14 compact, efficient and scalable representation and present a fast algorithm for their computation. On examples of  
15 several model organisms and two bacterial pan-genomes, we show that, compared to the best existing  
16 representation, simplitigs provide a substantial improvement in the cumulative sequence length and their number,  
17 especially for graphs with many branching nodes. We demonstrate that this improvement is amplified with more  
18 data available. Combined with the commonly used Burrows-Wheeler Transform index of genomic sequences,  
19 simplitigs substantially reduce both memory and index loading and query times, as illustrated with large-scale  
20 examples of GenBank bacterial pan-genomes.

## 21 Keywords

22 Sequence analysis, de Bruijn graphs,  $k$ -mers, representation, scaling

## 23 Background

24 Advances in DNA sequencing started the golden age of biology when previously unobservable phenomena can be  
25 studied on an unprecedented scale. However, sequencing capacity has been growing faster than computer  
26 performance and memory, and also faster than available human resources. Today large amounts of sequence data  
27 are available. Consequently, traditional sequence-based representations and sequence alignment-based techniques  
28 [1–3] have become less suitable for real-life scenarios due to the space- and time-complexities they impose and  
29 their inefficiency in handling polymorphisms.

30

31 De Bruijn graphs provide an elegant solution for genomic data representation. They build on top of the concept of  $k$   
32 -mers, which are substrings of a fixed length  $k$  of the genomic strings to be represented, such as sequencing reads,  
33 genomes, or transcriptomes. For a given  $k$ -mer set, the corresponding de Bruijn graph is a directed graph with the  
34  $k$ -mers being vertices and  $k - 1$  long overlaps between pairs of these  $k$ -mers indicating edges (Methods). There is  
35 an obvious correspondence between  $k$ -mer sets and de Bruijn graphs, and we can use both terms interchangeably.  
36 If  $k$  is chosen appropriately, de Bruijn graphs capture substantial information about the sequenced molecules as  
37 these correspond to some walks in the graph.

38

39 The use of de Bruijn graphs is ubiquitous in sequence analysis. Genome assembly uses the property that sequenced  
40 molecules form paths [4–6], which is exploited in numerous modern assemblers [7–12]. On the other hand,  
41 alignment-free sequence comparison follows the idea that similar sequences share common  $k$ -mers, and  
42 comparing de Bruijn graphs provides thus a good measure of sequence similarity [13]. This involves applications of  
43 de Bruijn graphs to variant calling and genotyping [14–18], transcript abundance estimation [19], and metagenomic  
44 classification [20–23]. In the latter application,  $k$ -mer-based classifiers perform best among all classifiers in  
45 inferring abundance profiles [24], which suggests that de Bruijn graphs truthfully approximate the graph structures  
46 of bacterial pan-genomes, even if constructed from noisy assemblies from incomplete databases. Even if more  
47 advanced pan-genome graph representations are available, such as variation graphs [25], de Bruijn graphs with  
48 large  $k$ -mer lengths are still useful for indexing [26,27].

49

50 As de Bruijn graphs are one of the primary data structures in much of sequence analysis, the efficiency of many  
51 algorithms is directly tied to the efficiency of computation and representation of the graph. De Bruijn graphs can be  
52 readily computed through a scan of the datasets including the raw reads, genomes or multiple sequence files. In  
53 practice, such a scan often consists in  $k$ -mer counting as this allows efficient denoising of the graph, e.g., by  
54 removing low-frequency  $k$ -mers corresponding to sequencing errors in the reads. Algorithms for  $k$ -mer counting  
55 have been extensively studied and many well-engineered software solutions are available [28–36].

56

57 On the other hand, de Bruijn graph representations have received much less attention. The most commonly used  
58 representation are unitigs, which are strings resulting from compaction of  $k$ -mers along maximal paths with  
59 non-branching nodes [37,38]. Unitigs have many advantages: the representation is “textual”, in the form of a set of  
60 sequences that contain each  $k$ -mer exactly once while preserving graph topology. However, unitigs impose large  
61 resource overhead for many types of de Bruijn graphs and do not scale well when a lot of variation is included.  
62 Specifically, with a high proportion of branching nodes, unitigs become fragmented, in an extreme case up to the  
63 level of individual  $k$ -mers. Subsequently, unitig computation and storage may require inappropriately large  
64 resources and become prohibitive in variation-heavy applications, e.g., in bacterial pan-genomics.

65

66 In this paper, we propose simplitigs as a compact, efficient and scalable representation of de Bruijn graphs.  
67 Simplitigs generalize the unitig representation by relaxing the restriction of stopping at branching nodes. We  
68 present an algorithm for rapid simplitig computation from a  $k$ -mer set and implement it in a tool called ProphAsm.  
69 ProphAsm proceeds by loading a  $k$ -mer set into memory and a greedy enumeration of maximal vertex-disjoint  
70 paths in the associated de Bruijn graph. We use ProphAsm to evaluate the improvement of simplitigs over unitigs,  
71 in terms of two main characteristics: the cumulative sequence length (CL) and the number of sequences (NS). We  
72 demonstrate that greedily computed simplitigs are close to theoretical bounds in practical applications and provide,  
73 compared to unitigs, a substantial improvement in memory requirements and speed in applications such as  $k$ -mer  
74 matching.

## 75 Results

### 76 The concept of simplitigs

77 We developed the concept of simplitigs to efficiently represent de Bruijn graphs of sequence data (**Fig. 1**).

78 Simplitigs are a generalization of unitigs and correspond to spellings of vertex-disjoint paths covering a given de  
79 Bruijn graph (**Fig. 1a**, Methods). Consequently, maximal simplitigs are such simplitigs where no two simplitigs can  
80 be merged by a  $(k - 1)$  overlap (Methods). Note that unitigs and  $k$ -mers are also simplitigs, but not maximal in  
81 general (**Fig. 1b**). The main conceptual difference between maximal simplitigs and maximal unitigs is that  
82 simplitigs are not limited by branching nodes, which allows for further compaction, with a benefit increasing  
83 proportionally to the amount of branching nodes in the graph.

84

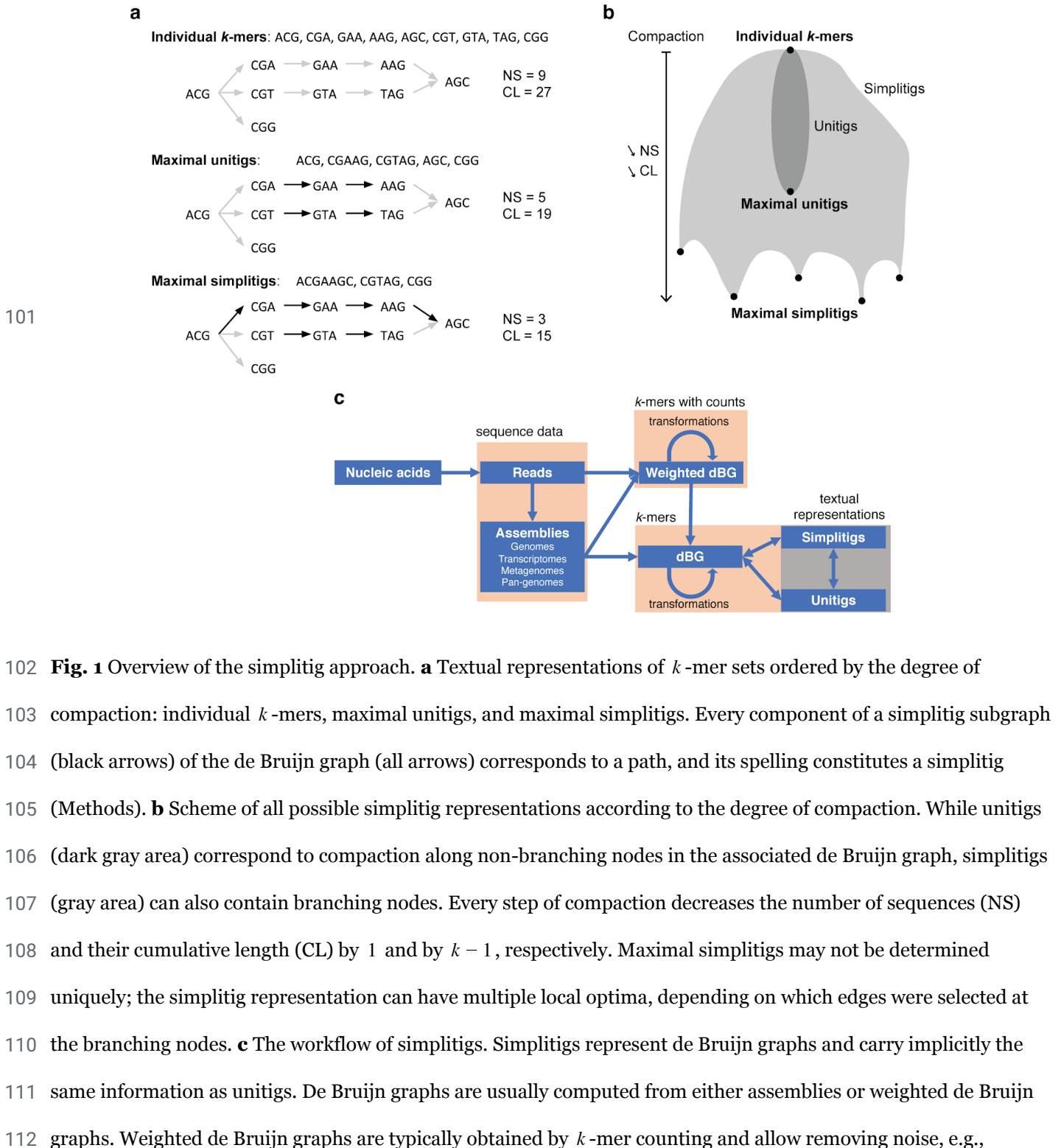
85 To compare simplitig and unitig representations, we created a benchmarking procedure based on the two  
86 characteristics: the number of sequences (NS) and their cumulative length (CL) (Methods, example in **Fig. 2**).

87 While NS determines the number of records to be kept in memory, CL largely determines the total memory needed.  
88 NS and CL are readily bounded from below by one and by the number of  $k$ -mers, respectively, and they are also  
89 tightly connected ((eq 1) in Methods). As every step of compaction decreases both NS and CL (**Fig. 1b**, Methods),  
90 we can optimize them jointly. However, finding an optimal simplitig representation translates to the vertex disjoint  
91 path coverage problem. While NP-hard for general graphs (by reduction from the well-known NP-hard problem of  
92 computing a Hamiltonian path), the problem may be tractable for observed de Bruijn graphs (Methods).

93

94 Since practical applications do not require optimal simplitigs, we prioritized speed and designed a greedy algorithm  
95 for their rapid computation (**Alg. 1**, Methods). In an iterative fashion, the algorithm selects an arbitrary  $k$ -mer as a  
96 seed of a new simplitig and keeps extending it forwards and then backwards as long as possible, while removing the  
97 already used  $k$ -mers from the set. This process is repeated until all  $k$ -mers are covered. Loading  $k$ -mers into  
98 memory and simplitig computation are linear in the length of the input and the number of  $k$ -mers, respectively,

99 and the memory footprint is linear in the number of  $k$ -mers. We implemented **Alg. 1** in a program called  
 100 ProphAsm, available at <https://github.com/prophyle/prophasm>.



113 low-frequency  $k$ -mers, which frequently originate in sequencing errors.

114

115 **Alg. 1** Greedy computation of maximal simplitigs for a  $k$ -mer set or a de Bruijn graph. In an iterative fashion, the  
116 algorithm draws a  $k$ -mer from the set of canonical  $k$ -mers  $K$ , uses it as a seed of a new simplitig, and then keeps  
117 extending the simplitig forwards as long as possible, and then backwards, while removing the already used  
118 canonical  $k$ -mers from  $K$ .

```
119 Function extend_simplitig_forwards (K, simplitig):
120     extending = True
121     while extending:
122         extending = False
123         q = suffix (simplitig, k-1),
124         for x in ['A', 'C', 'G', 'T']:
125             canon_kmer = canonical (q + x)
126             if canon_kmer in K:
127                 extending = True
128                 simplitig = simplitig + x
129                 K.remove (canon_kmer)
130             break
131     return K, simplitig
132
133 Function compute_maximal_simplitig_from_kmer (K, seeding_kmer):
134     simplitig = initial_kmer
135     K, simplitig = extend_simplitig_forwards (K, simplitig)
136     simplitig = reverse_completent (simplitig)
137     K, simplitig = extend_simplitig_forwards (K, simplitig)
138     return K, simplitig
139
140 Function compute_simplitigs (input_kmers):
141     K = {}
142     for kmer in input_kmers:
143         K.add (canonical (kmer))
144     maximal_simplitigs = {}
145     while |K| > 0:
146         seeding_kmer = K.pop ()
147         K, simplitig = compute_maximal_simplitig_from_kmer (K, seeding_kmer)
148         maximal_simplitigs.add (simplitig)
149     return maximal_simplitigs
```

## 150 **Simplitigs of model organisms**

151 We first evaluated simplitig and unitig representations on assemblies of six model organisms (**Fig. 2**). As different  
152 applications of de Bruijn graphs call for different  $k$ -mer lengths, we sought to characterize the NS and CL scaling  
153 for both representations with  $k$  growing, as well as the effect of the species' genome size. Therefore, selected model  
154 organisms were evaluated in an increasing order of the genome size and benchmarked for both representations on a  
155 range of  $k$ -mer lengths corresponding to common alignment-free-based applications [19,20,39].

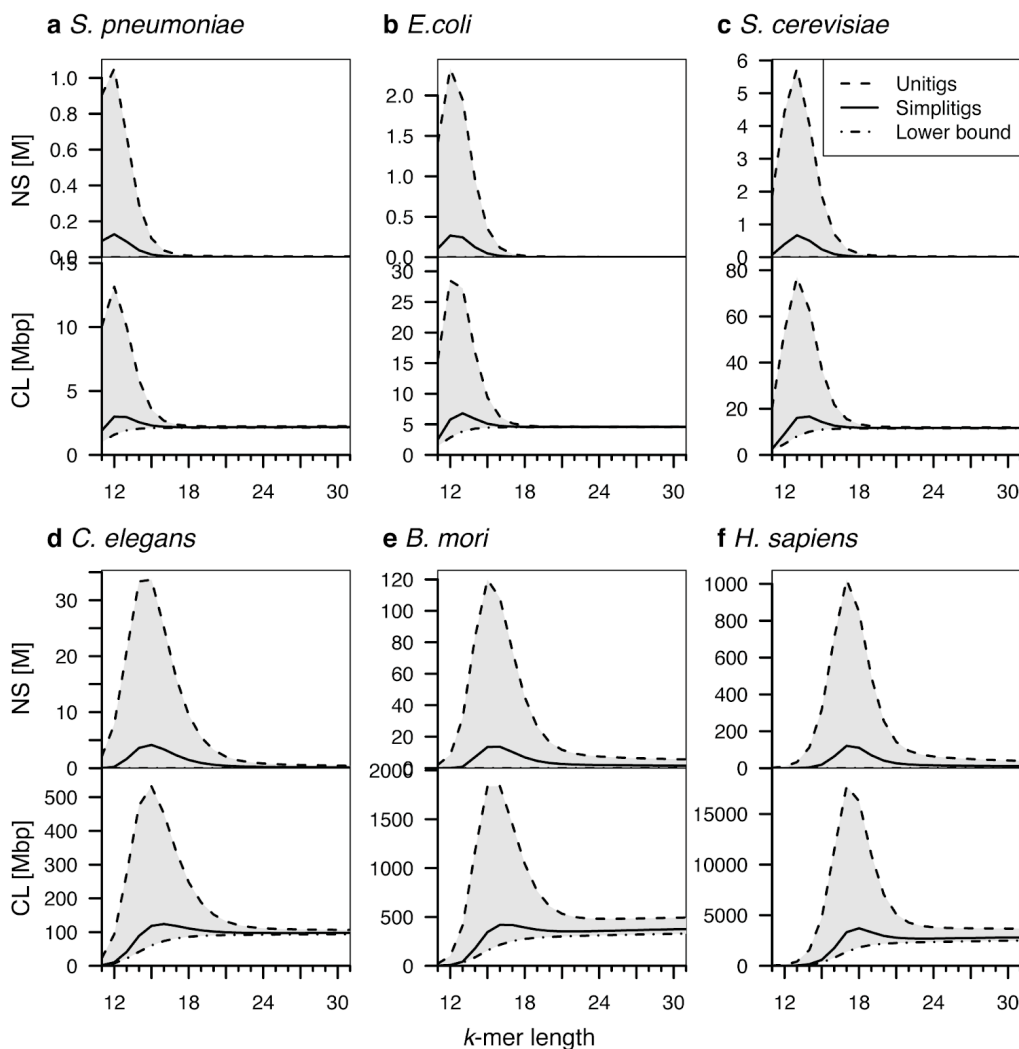
156

157 We observe that simplitigs always provide substantially better performance than unitigs (**Fig. 2**). In particular, they  
158 quickly approach the theoretical lower bounds for both characteristics tested. Every data set has a range of  $k$ -mer  
159 lengths where the difference between simplitigs and unitigs is very large, and after a certain threshold, the  
160 difference almost vanishes. While for short genomes this threshold is located at smaller  $k$ -mer lengths than those  
161 typically used in alignment-free applications (e.g.,  $k \approx 17$  for *E. coli*), for bigger genomes this threshold has not  
162 been attained on the tested range and seems to be substantially shifted towards large  $k$ -mers (e.g., *B. mori*).

163

164 Interestingly, maxima of the NS and CL values for both representations occur very close to the value  $k = \log_4 G$ ,  
165 where  $G$  is the genome size (**Fig. 2**). This is readily explained by edge saturation: for values of  $k$  up to  $\log_4 G$ , an  
166 overwhelming fraction of all  $4^k$   $k$ -mers belong to the genome, which makes the de Bruijn graph branch at nearly  
167 every node. As a consequence, unitigs are then essentially reduced to individual  $k$ -mers and their number grows  
168 exponentially whereas simplitigs stay compact on the whole range of  $k$ -mer lengths. Starting from  $k = \log_4 G$  the  
169 graph starts to form longer non-branching paths, which drives down the NS and CL of unitigs, and they approach  
170 those of simplitigs. However, the difference between simplitigs and unitigs in their count and length may stay  
171 considerable even for larger values of  $k$ , especially in case of large eukaryotic genomes.





172

173 **Fig. 2** Comparison of the simplitig and unitig representations for selected model organisms and a range of  $k$ -mers.  
174 The number of sequences (NS, millions) and their cumulative length (CL, megabase pairs) for both representations  
175 of six model organisms ordered by their genome size: **a** *Streptococcus pneumoniae*, 2.22 Mbp; **b** *Escherichia coli*,  
176 4.64 Mbp; **c** *Saccharomyces cerevisiae*, 12.2 Mbp; **d** *Caenorhabditis elegans*, 100 Mbp; **e** *Bombyx mori*, 482 Mbp;  
177 and **f** *Homo sapiens*, 3.21 Gbp. The CL lower bound corresponds to the number of  $k$ -mers. Full results are available  
178 in **Additional File 1**.

## 179 Performance assessment

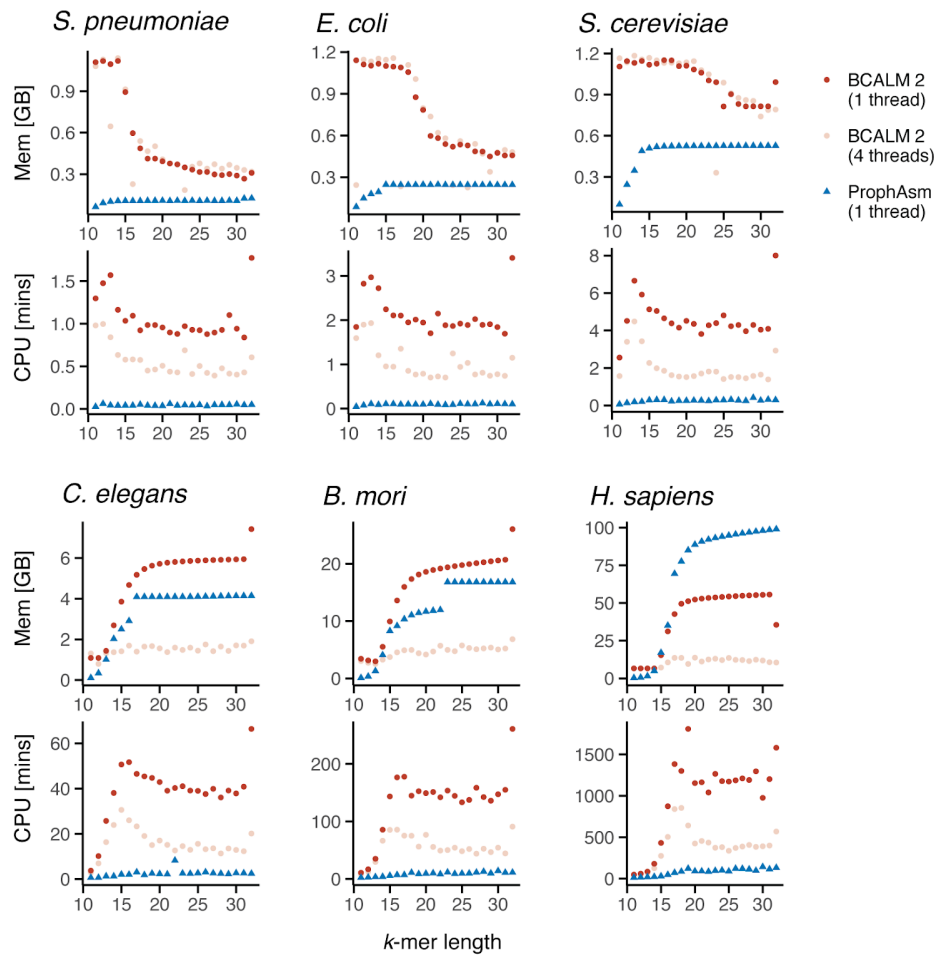
180 We then analyzed the computational resources that had been used to compute the simplitigs using ProphAsm and  
181 the unitigs using BCALM (**Fig. 3**). Both representations were computed in the environment of a computational  
182 cluster, with individual experiments deployed as parallel jobs using SLURM (Methods). Even though we primarily  
183 focused on the total CPU time, we also tested BCALM using 4 threads to evaluate the effect of parallelization.

184  
185 We observe that throughout our experiments ProphAsm outcompeted BCALM, both in terms of the memory  
186 consumption and the CPU time, even if BCALM was run with 4 threads (**Fig. 3**). The only exception was the  
187 memory consumption for *H. sapiens*, suggesting that BCALM may be more memory-efficient for large genomes.  
188 Importantly, ProphAsm used a consistent amount of resources: memory 38–51 bytes of RAM per  $k$ -mer in the  
189 dataset, a limited CPU time and no additional disk space. On the other hand, BCALM's resource usage was less  
190 predictable and less consistent across experiments and required frequent trial-and-error resource adjustments and  
191 re-running (Methods).

192  
193 Besides the comparatively high memory and CPU time requirements, the most challenging resource was the disk  
194 space. In order to fit within the available disk capacity, we reduced the BCALM consumption using the '-maxdisk'  
195 parameter. However, despite that we requested BCALM to use no more than 30 GB per experiment, a random  
196 manual inspection revealed a substantially higher use – for instance, 116 GB of disk space for *H. sapiens* with  
197  $k = 17$  and 4 threads. The extensive disk space consumption also disabled a similar comparison on a desktop  
198 computer.

199  
200 Overall, the better performance of ProphAsm can be explained by simplitigs being fundamentally easier to compute  
201 than unitigs and by BCALM being optimized for a particular class of de Bruijn graphs that emerge in assembly-like  
202 applications. As the ProphAsm resource usage is highly predictable, it appears to be more suitable for the use on  
203 clusters in many parallel instances, unless the number of  $k$ -mers in the dataset exceeds a critical threshold  
204 determined by the available RAM.

205



206

207 **Fig. 3** Comparison of CPU time and memory consumption of ProphAsm and BCALM. Resources to compute  
208 unitigs using BCALM (using one and four threads) and simplitigs using ProphAsm (using one thread) of the six  
209 model organisms. Full results are available in **Additional File 2**.

## 210 **Simplitigs of bacterial pan-genomes**

211 We then sought to evaluate the impact of additional variation in a de Bruijn graph (**Fig. 4**). Such variation may  
212 originate in polymorphisms, varying gene content in a population of genomes that are represented jointly, in  
213 haplotypes of viral quasispecies, or in sequencing errors in case of graphs constructed directly from sequencing  
214 reads. In all these cases, many nodes of the de Bruijn graph become branching and new paths emerge. To model  
215 gradually increasing variation, we used bacterial pan-genomes with different levels of sampling. Given the high  
216 diversity and variability of bacteria, de Bruijn graphs provide a convenient option for computational  
217 pan-genomes [40]. Such pan-genomes can be constructed from draft assemblies or even directly from sequencing  
218 reads, and thanks to bacterial genomes being short and haploid, the information captured by the graphs is sufficient  
219 for many analyses.

220

221 We first constructed a pan-genome of *N. gonorrhoeae*, and characterized unitigs and simplitigs as a function of  
222 pan-genome size (**Fig. 4a**). We used 1,102 draft assemblies of clinical isolates from the Prevention's Gonococcal  
223 Isolate Surveillance Project [41], from which we built a series of de Bruijn graphs using an increasing number of  
224 genomes. Consistent with previous experiments (**Fig 2ab**,  $k = 31$ ), both representations perform comparably well  
225 when only one bacterial genome is included (**Fig. 4a**). However, as the number of genomes or  $k$ -mers grows, the  
226 NS and CL grow as well, but with an increasing gap between unitigs and simplitigs; importantly, the latter stay close  
227 to the theoretical lower bounds. When the pan-genome size is measured via the number of genomes included, the  
228 CL and NS resemble logarithmic functions for both unitigs and simplitigs (**Fig. 4a**, left-hand column). However,  
229 when the number of  $k$ -mers included is used instead, the NS and CL functions act as affine functions (**Fig. 4a**,  
230 right-hand column). This suggests that a pan-genome  $k$ -mer count and a species-specific slope may be used as the  
231 predictors of simplitig performance in future applications.

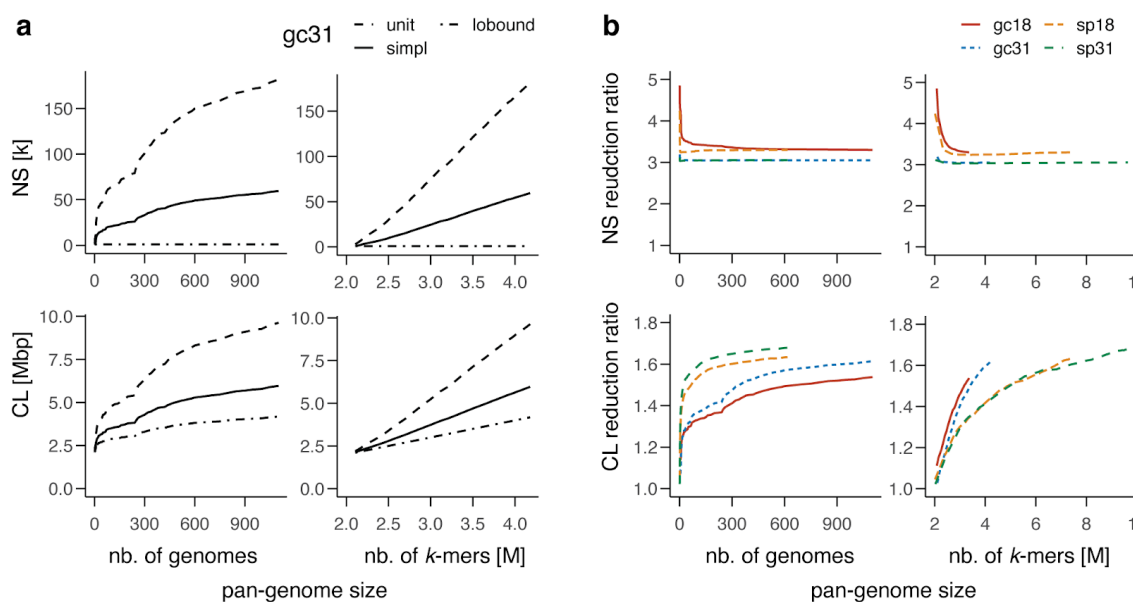
232

233 To analyze the relative benefit of simplitigs with growing de Bruijn graphs, we evaluated the NS and CL reduction  
234 ratio of simplitigs over unitigs in different configurations (**Fig. 4b**). We used the same *N. gonorrhoeae* dataset and  
235 considered also another dataset of *S. pneumoniae*, consisting of 616 draft Illumina assemblies of isolates from a

236 carriage study of children in Massachusetts, USA [42,43]. For both species and for  $k = 18, 31$ , we constructed a  
 237 series of de Bruijn graphs as previously, but this time we visualized the NS and CL reduction ratios. In all cases, the  
 238 NS reduction ratios eventually stabilized at values close to 3, following an L-shape ( $k = 18$ ) or being almost constant  
 239 ( $k = 31$ ). The CL reduction ratio admitted approximately a logarithmic dependence on the number of genomes and  
 240 still resembled a linear dependence on the number of  $k$ -mers. Overall, these experiments provided further evidence  
 241 that the benefit of simplitigs over unitigs grows with the increased proportion of branching nodes in a de Bruijn  
 242 graph or with increasing data in case of pan-genome reference structures.

243

244



245 **Fig. 4** Scaling of simplitigs and unitigs of bacterial pan-genomes as the pan-genome size grows with a better  
 246 sampling and more within-species variation. **a** Number of sequences (NS, thousands) and their cumulative length  
 247 (CL, megabase pairs) for simplitigs (simpl) and unitigs (unit) and the lower bound (lobound) of *N. gonorrhoeae*  
 248 and  $k = 31$ , as a function of the number of genomes (left) and  $k$ -mers (right, millions) included. **b** Reduction ratio  
 249 of simplitigs over unitigs for *S. pneumoniae* (sp) and *N. gonorrhoeae* (gc) and  $k = 18, 31$ , as a function of the  
 250 number of genomes (left) and  $k$ -mers (right, millions) included. Full results are available in **Additional File 3**.

## 251 Application of simplitigs for $k$ -mer search

252 Finally, we sought to demonstrate the benefit of simplitigs in a real application. A major use of de Bruijn graphs  
253 consists in  $k$ -mer matching, which requires the graphs to act as a membership data structure. As both simplitigs  
254 and unitigs are text-based representations,  $k$ -mer queries can be implemented using an arbitrary full-text  
255 index [44], notably a Burrows-Wheeler Transform index [45] (sometimes referred to as an FM-index). Here, we  
256 used the index of BWA [46], as one of the best-engineered solutions available, to analyze the impact of unitig  
257 replacement by simplitigs.

## 258 Single pan-genome

259 We first evaluated the simplitig improvement on the same *N. gonorrhoeae* pan-genome (**Fig. 5**). We considered  
260 four different  $k$ -mer lengths  $k = 19, 23, 27, 31$  and for each of them, we built three pan-genome representations from  
261 the original draft genome assemblies: first, we merged the assemblies as the most straightforward approach to  
262 collect all  $k$ -mers; then we computed pan-genome representations by unitigs and simplitigs. As all the three  
263 representations carry the same  $k$ -mer set, a full-text index built upon them should provide the same results, but  
264 with a performance reflecting the differences in NS and CL.

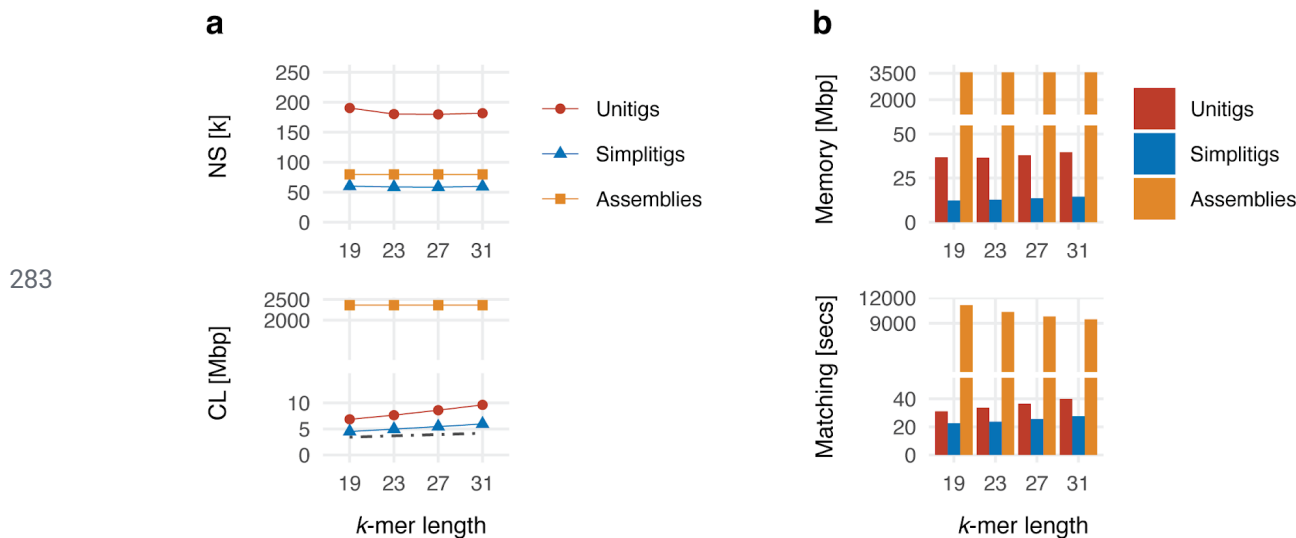
265

266 We analysed the NS and CL characteristics of the computed representations (**Fig. 5a**). Simplitigs improved NS and  
267 CL over unitigs by factors of  $3.1\times$ – $3.2\times$  and  $1.5\times$ – $1.6\times$ , respectively, consistent with **Fig. 4**. This suggests that  
268 memory required for unitigs would be approximately by 50% higher than for unitigs across all the values of  $k$   
269 considered. We also compared simplitigs and unitigs to assemblies; both improved the CL by two orders of  
270 magnitude while the NS stayed comparable for simplitigs and increased twofold for unitigs.

271

272 We then evaluated  $k$ -mer query performance and memory footprint (**Fig. 5b**). For each representation, we  
273 constructed a standard BWA index and matched 10 million random  $k$ -mers from the pan-genome using BWA  
274 fastmap [47], with no restriction on maximum size of the suffix-array interval to ensure evaluation correctness  
275 (Methods). Surprisingly, simplitigs improved memory by factors of  $2.7\times$  –  $3.0\times$  (**Fig. 5b**), thus twice as much as we

276 previously anticipated. This is explained by the fact that the underlying full-text engine has to keep information  
277 about individual sequences in memory as separate records. As NS grows, it has a negative impact on both the  
278 memory footprint and query speed. Nevertheless, since simplitigs provided a substantial reduction in NS over  
279 unitigs, this overhead has been reduced. Therefore, the excessive number of unitigs observed throughout our  
280 experiments (**Fig. 1** and **Fig. 2**) provides a further argument for replacing unitigs by simplitigs when possible.  
281  
282



284 **Fig. 5** *K*-mer queries for the *N. gonorrhoeae* pan-genome on top of the draft assemblies, unitigs, and simplitigs.

285 **a** Characteristics of the obtained unitigs and simplitigs: number of sequences (NS, thousands) and their cumulative  
286 length (CL, megabase pairs). The dot-dash line depicts the CL lower bound corresponding to the number of *k*-mers.

287 **b** Time and memory footprint of BWA *k*-mer queries (10 million *k*-mers). Full results including relative

288 improvements are available in **Additional File 4**.

## 289 Multiple pan-genomes

290 We evaluated the performance of the simplitig representation for simultaneous indexing of a large number bacterial  
291 pan-genomes (**Fig. 6**). We downloaded all complete bacterial genomes from Genbank that had not been excluded  
292 from RefSeq (as of May 2020; 9,869 records out of which 9,032 had genomic sequences available; Methods). We  
293 restricted ourselves to complete genomes as draft genomes in Genbank are substantially impacted by false genetic  
294 variability [48–50] that is particularly common in bacterial studies, mainly due to the contaminant DNA [51]. By  
295 grouping individual genomes per species, we obtained 3,179 bacterial pan-genomes which we call the “All” dataset.  
296 After computing simplitigs and unitigs per species, we merged the obtained representations and constructed  
297 indexes using BWA; all this was done for  $k = 19, 23, 27, 31$  to evaluate the impact of the  $k$ -mer length. As the unitig  
298 index could not fit into RAM of our desktop computer for any  $k$ , we also created the “Solid” dataset by omitting  
299 pan-genomes with less than 11 genomes; this resulted in 112 pan-genomes with 3,958 genomes. We provide all the  
300 constructed pan-genomes in the form of simplitigs on Zenodo ([10.5281/zenodo.3800713](https://zenodo.org/record/3800713)).

301

302 First we analyzed the obtained simplitig and unitig representations of both datasets (**Fig. 6a**). We observe that  
303 simplitigs provide substantial improvement in both test characteristics. In the Solid dataset, NS and CL were  
304 reduced by simplitigs by factors of 3.1–4.5 and 1.4–1.9, respectively; and in the All dataset, NS and CL were reduced  
305 by factors of 3.0–4.3 and 1.2–1.4, respectively; all consistent with the scaling observed previously (**Fig. 4, Fig. 5**).  
306 While the improvement in NS was almost identical in both datasets (consistent with the top-right graph in **Fig 4b**),  
307 the improvement in CL was clearly better in the Solid dataset. Indeed, as the vast majority of pan-genomes in the  
308 All dataset contained only one genome, the de Bruijn graphs had a comparatively low number of branching nodes,  
309 therefore the difference between simplitigs and unitigs was less striking (consistent with the values for small  
310 pan-genome sizes in **Fig 4b**). We also observe that, in contrast to unitigs,  $k$ -mer length had only little impact on  
311 the CL of simplitigs within the tested range, which provides better guarantees on required computational resources  
312 in future applications.

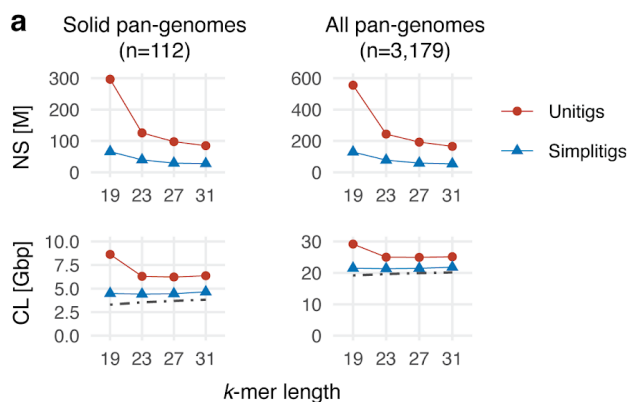
313

314 We then measured the performance of  $k$ -mer lookup (**Fig. 6b**). Both on a desktop and on a cluster, we evaluated  
315 memory footprints, index loading time and time to match ten million random  $k$ -mers from the index using BWA

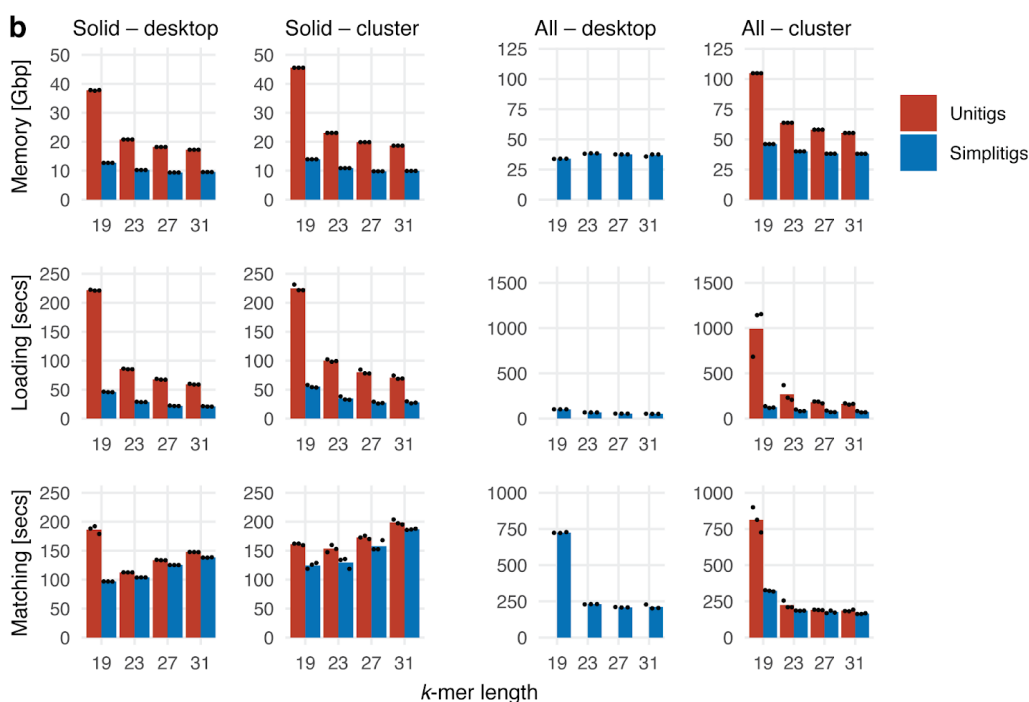


316 (Methods). We observed that simplitigs substantially improved the memory footprint and index loading times. For  
317  $k = 19$ , simplitigs largely improved the matching times, where the difference is caused by “ghost  $k$ -mers” on unitig  
318 borders; these were more common in this experiment due to the short  $k$ -mer length and the high number of  
319 unitigs. For higher  $k$ -mer lengths, simplitigs still provided a moderate improvement in the matching rate. We note  
320 that the query time with BWT-based  $k$ -mer indexes is dominated by high-frequency  $k$ -mers. As these are equally  
321 frequent in simplitigs and unitigs, the observed performance is similar unless many ghost  $k$ -mers emerge on  
322 sequence borders as seen previously. On the desktop, the unitig All-index could not be evaluated as it did not fit into  
323 memory, and the outlier for  $k = 19$  may be the result of memory swapping.

324



325



326 **Fig. 6** *K*-mer queries for multiple pan-genomes indexed simultaneously. Bacterial pan-genomes were computed  
 327 from the complete Genbank assemblies per individual species. While the All dataset comprises all pan-genomes  
 328 with no restriction on their size, the Solid dataset comprises only those that contain at least 11 genomes.  
 329 **a** Characteristics of the obtained unitigs and simplitigs: number of sequences (NS, millions) and their cumulative  
 330 length (CL, gigabase pairs). The dot-dash line depicts the lower bound corresponding to the number of *k*-mers.  
 331 **b** Memory footprint, time of index loading and time of matching 10 million *k*-mers using BWA. The bars  
 332 correspond to the mean of three measurements (black dots). Full results including relative improvements are  
 333 available in **Additional File 5**.

## 334 Discussion

335 We introduced the concept of simplitigs, a generalization of unitigs, and demonstrated that simplitigs constitute a  
336 compact, efficient and scalable representation of de Bruijn graphs for various types of genomic datasets. The two  
337 representations share many similarities: they are text-based and individual strings correspond to spellings of  
338 vertex-disjoint paths. Both representations can be seen as irreversible transforms, taking a set of input strings and  
339 producing a new set of strings preserving the  $k$ -mer sets. In both cases, the resulting files can easily be manipulated  
340 using common Unix tools, compressed using standard compression techniques, and indexed using full-text indexes.  
341 The main difference consists in that simplitigs do not explicitly carry information about the topology of the de  
342 Bruijn graph. Also simplitigs are not expected to have direct biological significance – neighboring segments of the  
343 same simplitig may correspond to distant parts of the same nucleic acid or even to different ones. Nevertheless,  
344 unitigs can always be recomputed from simplitigs, but this step is not required for many common applications.  
345 Furthermore, a concept analogous to simplitigs, called disjointigs, was recently introduced in the context of genome  
346 assembly using A-Bruijn graphs [52,53], suggesting that simplitigs may be useful beyond the context of  
347 topology-oblivious applications.

348

349 We provided ProphAsm, a tool implementing a greedy heuristic to efficiently compute maximal simplitigs from a  $k$   
350 -mer set. ProphAsm is a spin-off of the ProPhyle software (<https://prophyle.github.io>, [22,54]) for metagenomic  
351 classification, allowing efficient indexing of  $k$ -mers propagated to individual nodes of a phylogenetic tree.  
352 ProphAsm presents a “naive” implementation of the greedy heuristic (**Alg 1**) that can be further improved. For  
353 instance, a hash table with better memory management may reduce the memory requirement by a factor of 2.5 [55].  
354 Additional memory reduction may be also achieved in a similar fashion as previously done for unitigs [38,56,57].  
355 Nevertheless, on the studied data, ProphAsm strikingly outcompeted BCALM in all characteristics measured, with  
356 the only exception of memory in the case of *H. sapiens*.

357

358 The observed lower requirements of ProphAsm can be attributed to two major differences between unitigs and  
359 simplitigs. First, unitigs, and BCALM as one of the reference programs, are designed for assembly applications,

360 where only a small proportion of nodes is usually branching. Moreover, experiments are usually run on a cluster in  
361 series, with multiple threads per experiment, rather than in parallel, and the tools can use extensive disk space.  
362 However, when unitigs are tested as a general representation with restricted disk space, the required resources  
363 grow substantially. Second, simplitigs are likely to be, from a computational perspective, fundamentally easier than  
364 unitigs.

365

366 A challenging but also promising feature of simplitig representation is the ambiguity of maximal simplitigs. This is  
367 in sharp contrast to maximal unitigs, which are uniquely defined (up to the order, reverse complementing, and  
368 cycles). In practice, every algorithm for simplitig computation has to decide which edge will be included at each  
369 branching node. Here, we prioritized speed, simplitigs were constructed progressively and lexicographically  
370 minimal edges were used in a case of ambiguity. Therefore, the final maximal simplitigs were only dependent on the  
371 choice of seeding  $k$ -mers; these are determined by the specific implementation of ‘std::unordered\_set’ in the C++  
372 standard library. Nevertheless, characteristics other than speed could readily be prioritized instead. For instance, a  
373 more sophisticated heuristic may shift the CL and NS closer to the optimum. One could also aim at adding  
374 biological meaning to simplitigs, e.g., by preferring those paths that are better supported by sequence data. On the  
375 other hand, streaming algorithms for operations such as merging or intersecting may require specific prescribed  
376 forms of unitigs. Finally, simplitigs may be optimized for entropy in order to maximize their compressibility.

377

378 The data presented in this paper highlight the scaling of computational resources as more sequencing data become  
379 available. The studied *N. gonorrhoeae* dataset constitutes a relatively complete image of a bacterial population in a  
380 geographical region at a given time scale. As such, it can be used to model the “state of completion” of  $k$ -mer-based  
381 pan-genome representations. On the other hand, the multiple pan-genomes experiment shows how simplitigs  
382 perform when a large number of such pan-genomes, although in different states of completion, are considered  
383 simultaneously and queried using a single BWT index. Overall, the presented experiments allow us to predict the  
384 resources for species for which only limited sequence data are available at present, but more are likely to be  
385 generated in the future. Importantly, with more data available, comparative benefits of simplitigs over unitigs grow  
386 as we have shown throughout the paper. As the growth of public databases negatively impacts the accuracy of those

387 algorithms that operate on top of lossily represented de Bruijn graphs [58], simplitigs provide a promising solution  
388 offering both exactness and scalability.

389

390 In modern bioinformatics applications multiple de Bruijn graphs are often considered simultaneously; the resulting  
391 structure is usually referred to as a colored de Bruijn graph [14] and the associated data structures have been widely  
392 studied [59–70]. Although we touched upon this setting in the Multiple pan-genomes section, exploiting the  
393 similarity between individual de Bruijn graphs for further compression in simplitig-based approaches is to be  
394 addressed in future work. In many applications, including some of the traditional alignment-free methods [13,71], it  
395 is also desirable to consider  $k$ -mers with counts. In the context of de Bruijn graphs, this leads to the so-called  
396 weighted variant of the problem [72]. The fact that frequencies of overlapping  $k$ -mers are usually similar suggests  
397 that  $k$ -mers can be grouped based on frequencies and simplitigs constructed per group .

398

399 Independently and simultaneously with the work we present here, the simplitig representation was recently studied  
400 in [73] under the name “spectrum-preserving string sets”. The associated UST tool follows a similar greedy strategy  
401 to ProphAsm, although operating on unitigs constructed by BCALM rather than on the original  $k$ -mers. As we  
402 demonstrated throughout this paper, unitigs are prohibitive for highly branching de Bruijn graphs, where a  
403 simplitig construction through unitigs may create a burden on resources and easily become intractable. The paper  
404 presents a tighter lower bound on the cumulative length of the representation (CL), termed weight, taking into  
405 account the graph topology but requiring a computational overhead. The authors also studied the compression  
406 properties of simplitigs when combined with standard compression algorithms. On the other hand, the paper does  
407 not study the number of sequences (NS) and continuous scaling for parameters such as  $k$ -mer length or amount of  
408 additional variation included.

## 409 Conclusions

410 In this paper, we addressed the question of efficient and scalable representation of de Bruijn graphs. We showed  
411 that the state-of-the-art unitig representation may require adequately large computational resources, especially  
412 when de Bruijn graphs contain many branching nodes. We introduced simplitigs, which provide a more compact  
413 replacement in applications that do not require explicit information on the graph topology, such as alignment-free  
414 sequence comparison and  $k$ -mer indexing. We introduced a heuristic simplitig computation and showed on the  
415 examples of model species that unless the genome is large, even a naive implementation outperforms BCALM, the  
416 main state-of-the-art tool for unitigs. We then studied applications to bacterial pan-genomics and showed that the  
417 utility of simplitigs compared to unitigs grows as more data are available. Finally, we demonstrated on the example  
418 of full-text  $k$ -mer indexing that simplitigs can substantially reduce computational resources and allow  
419 computations in situations where unitigs would bring unaffordable costs when exactness should be preserved.

420

421 Our work opens many questions and future directions. The presented algorithm for simplitig computation can be  
422 improved, parallelized, de-randomized to ensure reproducibility. We anticipate more theoretical advances in the  
423 analysis of the minimum vertex-disjoint path cover problem and better mapping to results from other disciplines  
424 such as network sciences. We also anticipate improvements in the heuristic approaches that could simplify and  
425 parallelize simplitig computation. The nature of the algorithm suggests that simplitigs might be computed online  
426 directly from a stream of data such as sequencing reads. We anticipate better implementations and libraries for  
427 simplitigs that can be plugged into standard bioinformatics libraries for various programming languages. Another  
428 series of questions is related to low-memory transformations of computed simplitigs that would allow  
429 precomputing simplitigs on computer clusters, and tailoring to specific applications on standard computers; this  
430 includes decreasing  $k$ , performing set operations on top of simplitig sets and computing maximal unitigs from  
431 simplitigs. A substantial body of work can be anticipated in the direction of text indexing – we showed that  
432 simplitigs can be combined with full-text indexes; however, specialized simplitig indexes exploiting simplitig  
433 characteristics are yet to be developed. On the theoretical side, many perspectives are open in the direction of  
434 orthogonal algorithmic techniques, such as sketching [74,75], and in relation to other stringology concepts, such as

435 minimal absent words [76] and shortest superstrings [77]. Finally, simplitigs can serve as components in design of  
436 various specialized data structures; these can involve membership queries of classes of de Bruijn graphs and colored  
437 de Bruijn graphs, where simplitigs may encode not only the  $k$ -mers themselves, but also additional metadata such  
438 as  $k$ -mer frequencies or colors. Furthermore, simplitigs will facilitate new full-text-based data structures for  
439 approximate matching, based on the inclusion of the proximity variation, that would be completely intractable with  
440 unitigs. Overall, we anticipate that the simplitig representation will become a generic compact representation of de  
441 Bruijn graphs, in particular, in the context of large-scale sequence data search engines [69] and sequence data  
442 repositories such as those of NCBI and EBI.

## 443 Methods

### 444 De Bruijn graphs

445 All strings are assumed to be over the alphabet  $\{A, C, G, T\}$ . A  $k$ -mer is a string of length  $k$ . For a string  $s = s_1 \dots s_n$ ,  
446 we define  $pref_k(s) = s_1 \dots s_k$  and  $suf_k(s) = s_{n-k+1} \dots s_n$ . For two strings  $s$  and  $t$  of length at least  $k$ , we define the  
447 binary connectivity relation  $s \rightarrow_k t$  if and only if  $suf_k(s) = pref_k(t)$ . If  $s \rightarrow_k t$ , we define the  $k$ -merging operation  $\odot$  as  
448  $s \odot^k t = s \cdot suf_{|t|-k}(t)$ .

449

450 Given a set  $K$  of  $k$ -mers, the *de Bruijn graph* of  $K$  is the directed graph  $G = (V, E)$  with  $V = K$  and  
451  $E = \{(u, v) \in K^2 \mid u \rightarrow_{k-1} v\}$ . For every path  $p = (v_1, \dots, v_p)$  in  $G$ , the string  $v_1 \odot^{k-1} v_2 \odot^{k-1} \dots \odot^{k-1} v_p$  is called a *spelling* of  
452  $p$ . This definition of de Bruijn graphs is *node-centric*, as nodes are identified with  $k$ -mers and edges are implicit.  
453 Therefore, we can use the terms “ $k$ -mer set” and “de Bruijn graph” interchangeably.

### 454 Simplitigs

455 Consider a set  $K$  of  $k$ -mers and the corresponding de Bruijn graph  $G = (K, E)$ . A *simplitig graph*  $G' = (K, E')$  is a  
456 spanning subgraph of  $G$  that is acyclic and the in-degree and out-degree of any node is at most one. It follows from  
457 this definition that a simplitig graph is a vertex-disjoint union of paths, whose spellings we call *simplitigs*. A  
458 simplitig is called *maximal* if it cannot be extended forward or backward without breaking the definition of  
459 simplitig graph. In more detail, a simplitig  $u_1 \rightarrow_{k-1} u_2 \rightarrow_{k-1} \dots \rightarrow_{k-1} u_n$  is maximal if the following conditions hold

- 460
- 461 • either  $u_1$  has no incoming edges in  $G$ , or for any edge  $(v, u_1) \in E$ ,  $v$  belongs to another simplitig and it is  
462 not its last vertex,
  - 463 • either  $u_n$  has no outgoing edges in  $G$ , or for any edge  $(u_n, v) \in E$ ,  $v$  belongs to another simplitig and it is  
464 not its first vertex.

464 A *unitig* is a simplitig  $u_1 \rightarrow_{k-1} u_2 \rightarrow_{k-1} \dots \rightarrow_{k-1} u_n$  such that each of the nodes  $u_2, \dots, u_n$  has in-degree 1 and each of the  
465 nodes  $u_1, \dots, u_{n-1}$  has out-degree 1 in graph  $G$ . A maximal unitig is defined similarly.



## 466 Comparing simplitig and unitig representations

467 Simplitigs and unitigs representations were compared in terms of the number of sequences produced (NS) and their  
468 cumulative length (CL). For any set of simplitigs (i.e., not necessarily maximal ones), NS is bounded by 1 and  
469  $\#kmers$ , CL is bounded by  $\#kmers$  and  $k \cdot \#kmers$ . The upper bound corresponds to the state of maximal  
470 fragmentation, where every  $k$ -mer forms a simplig. The lower bound corresponds to the maximum possible degree  
471 of compaction, i.e., a single simplitig containing all  $k$ -mers.

472

473 NS and CL are readily connected by the following formula:

$$474 \quad CL = \#kmers + (k - 1) \cdot NS \quad (\text{eq 1})$$

475 As an important consequence, both characteristics are optimized simultaneously.

## 476 Greedy computation of simplitigs

477 The problem of computing maximal simplitigs that are optimal in CL (i.e., also in NS) corresponds to the minimum  
478 vertex-disjoint path cover problem [78]. This is known to be NP-hard in the general case, reducing from the  
479 Hamiltonian path problem. However, the complexity for de Bruijn graphs remains an open question. A greedy  
480 heuristic to compute maximal simplitigs has been used throughout this paper (**Alg. 1**). Simplitigs are constructed  
481 iteratively, starting from (arbitrary) seeding  $k$ -mers and being extended greedily forwards and backwards as long  
482 as possible.

## 483 ProphAsm implementation

484 ProphAsm is written in C++ and implements the greedy approach described above (**Alg. 1**).  $K$ -mers are encoded  
485 using `uint64_t` and stored in an `std::unordered_map`. The choice of extension nucleotides on branching nodes is  
486 done based on the lexicographic order. Therefore, the only source of randomness is the choice of seeding  $k$ -mers by  
487 `std::unordered_set::begin`; the C++ standard library makes no guarantees on which specific element is considered  
488 the first element. ProphAsm does not require any disk space to store intermediate data and its memory

489 consumption corresponded to 38–51 bytes per a unique  $k$ -mer (in dependence on the allocation), consistent with  
490 [55].

## 491 **Uni-directed and bi-directed models**

492 The uni-directed model, as presented above, is useful for introducing the concepts of unitigs and simplitigs, but is  
493 not directly applicable to data obtained using sequencing: since DNA is double-stranded, every string may come  
494 from either strand. At the level of  $k$ -mers, double-strandedness can be accounted for by using canonical  $k$ -mers,  
495 i.e., by pairing-up every  $k$ -mer with its reverse complement, typically done by taking the lexicographical minimum  
496 of the  $k$ -mer and its reverse complement. This subsequently requires redefining de Bruijn graphs to bi-directed  
497 de Bruijn graphs [79], which requires a more complex formalism.

## 498 **Correctness evaluation**

499 The correctness of simplitig computation can be verified using an arbitrary  $k$ -mer counter. Simplitigs have been  
500 computed correctly if and only if every  $k$ -mer is present exactly once and the number of distinct  $k$ -mers is the same  
501 as in the original datasets. The correctness of ProphAsm outputs was verified using JellyFish 2 [29].

## 502 **Experimental evaluation - model organisms and performance**

503 Reference sequences for six selected model organisms were downloaded from RefSeq and UCSC Genome Browser:  
504 *S. pneumoniae* str. ATCC 700669 (accession: NC\_011900.1, length 2.22 Mbp), *E. coli* str. K-12 (accession:  
505 NC\_000913.3, length: 4.64 Mbp), *S. cerevisiae* (accession: NC\_001133.9, length: 12.2 Mbp), *C. elegans* (accession:  
506 GCF\_000002985.6, length: 100 Mbp), *B. mori* (accession: GCF\_000151625.1, length: 482 Mbp), and *H. sapiens*  
507 (HG38, <http://hgdownload.soe.ucsc.edu/goldenPath/hg38/bigZips/hg38.fa.gz>, length: 3.21 Gbp). For each  
508 genome, simplitigs and unitigs were computed using ProphAsm and BCALM, respectively, for a range of  $k$ -mer  
509 lengths [11,31].

510

511 Individual experiments were run in parallel on the Harvard Medical School O2 cluster using Snakemake [80] and  
512 SLURM. ProphAsm and BCALM were run with the following parameters, respectively: ‘-k {kmer-length}’ and

513 '-kmer-size {kmer-length} -abundance-min 1 -nb-cores {cores} -max-disk 30000'. As BCALM requires a large  
 514 undocumented amount of disk space, we used the -max-disk parameter to make a parallel execution of many  
 515 BCALM jobs feasible. The SLURM specifications of resource allocation for individual species were iteratively  
 516 adjusted until all jobs would finish; the final required resources are provided in **Supplementary Table 1**. Time  
 517 and memory consumption of jobs were measured independently using GNU Time. Individual jobs were deployed to  
 518 computational nodes with different hardware configurations, which are specified in **Additional File 2**.

519

520 **Supplementary Table 1** SLURM resource allocation for ProphAsm and BCALM2 for the performance evaluation.

Species	ProphAsm (1 core)		BCALM (1 core)		BCALM (4 cores)	
	mem [GB]	cpu [hours]	mem [GB]	cpu [hours]	mem [GB]	cpu [hours]
<i>S. pneumoniae</i>	10	1	10	1	10	1
<i>E. coli</i>	10	1	10	1	10	1
<i>S. cerevisiae</i>	10	1	10	1	10	1
<i>C. elegans</i>	10	1	10	2	10	2
<i>B. mori</i>	20	2	30	10	30	2
<i>H. sapiens</i>	120	4	100	48	100	18

521

## 522 Experimental evaluation - bacterial pan-genomes

523 First, 1,102 draft assemblies of *N. gonorrhoeae* clinical isolates (collected from 2000 to 2013 by the Centers for  
 524 Disease Control and Prevention's Gonococcal Isolate Surveillance Project [41], and sequenced using Illumina  
 525 HiSeq) were downloaded from Zenodo [81]. Second, 616 draft assemblies of *S. pneumoniae* isolates (collected from  
 526 2001 to 2007 for a carriage study of children in Massachusetts, USA [42,43], and sequenced using Illumina HiSeq)  
 527 were downloaded from the SRA FTP server using the accession codes provided in Table 1 in [43]. For each of these  
 528 datasets, an increasing number of genomes was being taken and merged, and simplitigs and unitigs computed using  
 529 ProphAsm and BCALM, respectively. This experiment was performed for  $k = 18$  and  $k = 31$ . To avoid excessive  
 530 resource usage the functions were evaluated at selected points in an increasing distance: for intervals [10, 100] and  
 531 [100,+∞] only multiples of 5 and 20 were evaluated, respectively.

## 532 Experimental evaluation - full-text $k$ -mer queries

533 In the single pan-genome experiment, the same 1,102 assemblies of *N. gonorrhoeae* were merged into a single file.  
534 ProphAsm and BCALM were then used to compute simplitigs and unitigs, respectively, from this file for  
535  $k = 19, 23, 27, 31$ . Each of the three obtained FASTA files (assemblies, simplitigs, and unitigs) was used to  
536 construct a BWA index, which was then queried for  $k$ -mers using 'bwa fastmap -l {kmer-length}'. We used a  
537 modified version of BWA fastmap that reports both the time of index loading and the time of querying  
538 (<http://github.com/karel-brinda/bwa>, commit e1f907c). Query  $k$ -mers were generated from the same pan-genome  
539 using WGSim (version 1.10, with the parameters '-h 0 -S 42 -r 0.0 -l {kmer-length} -N 10000000 -e 0').

540

541 For the multiple pan-genome experiment, a list of available bacterial assemblies was downloaded from  
542 [ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/bacteria/assembly\\_summary.txt](ftp://ftp.ncbi.nlm.nih.gov/genomes/genbank/bacteria/assembly_summary.txt) (2020/05/05). For all assemblies  
543 marked as complete (i.e., the "assembly\_level" column equal to "Complete genome") and present in RefSeq (i.e., an  
544 empty value in the column "excluded\_from\_refseq"), directory urls and species names were extracted (n=9,869).  
545 These were then used to download the genomes of the isolates using RSync, restricting to genomic sequences only  
546 (i.e., files matching '\*v?\_genomic.fna.gz', n=9,032). The downloaded assemblies were then merged per species in  
547 order to collect  $k$ -mers of individual pan-genomes and used for computing simplitigs and unitigs using ProphAsm  
548 and BCALM, respectively. The obtained simplitig and unitig files were then merged per categories (e.g., simplitigs  
549 for  $k=19$ ) and used to construct a BWA index. The obtained indexes were queried for 10 million  $k$ -mers using BWA  
550 fastmap as previously. The  $k$ -mers were generated from the original assemblies of randomly selected 100 genomes  
551 using DWGsim [82] (version 0.1.11, with the parameters '-R 0 -e 0 -r 0 -X 0 -y 0 -H -z 42 -m /dev/null -N  
552 10000000 -l {k} -2 0'); the randomization was performed using 'sort -R'.

553

## 554 Computational setup

555 The experiments were performed on the HMS O2 research high-performance cluster and on an iMac 4.2 GHz  
556 Quad-Core Intel Core i7 with 40 GB RAM. The reproducibility of computation was ensured using BioConda [83].

557 All benchmarking was performed using ProphAsm 0.1.1 (commit ea28b708) and BCALM 2.2.2 (commit febf79a3).

558 Time and memory footprint were measured using GNU Time.

## 559 **Declarations**

### 560 **Ethics approval and consent to participate**

561 Not applicable.

### 562 **Availability of data and materials**

#### 563 **Additional files**

564 **Additional File 1.** Detailed information for the single genome experiment: NS, CL and #kmers for unitigs and  
565 simplitigs as a function of  $k$  for the 6 species: **a** *S. pneumoniae*, **b** *E. coli*, **c** *S. cerevisiae*, **d** *C. elegans*, **e** *B. mori*,  
566 and **f** *H. sapiens*.

567 **Additional File 2.** Detailed information for the performance comparison. **a** CPU time and memory consumption  
568 (both measured by GNU Time and Snakemake) as a function of species, method, number of threads, and  $k$ -mer  
569 length, including the used computational node. **b** Hardware specifications for individual computational nodes.

570 **Additional File 3.** Detailed information for the pan-genome scaling experiment: **a** *N. gonorrhoeae*,  $k = 18$ ; **b** *N.*  
571 *gonorrhoeae*,  $k = 31$ ; **c** *S. pneumoniae*,  $k = 18$ ; **d** *S. pneumoniae*,  $k = 31$ .

572 **Additional File 4.** Detailed information for the single pan-genome  $k$ -mer indexing experiment. **a** Characteristics  
573 of the resulting simplitigs and unitigs for  $k = 19, 23, 27, 31$ . **b** Memory footprint, index loading time and time to  
574 query 10 million  $k$ -mers using BWA.

575 **Additional File 5.** Detailed information for the multiple pan-genomes  $k$ -mer indexing experiment. **a** List of all  
576 genomes used for building the pan-genomes (accession code, version, species, filename, number of sequences,  
577 genome size [bp]); **b** List of species and the number of genomes included. **c** Characteristics of the resulting  
578 simplitigs and unitigs of individual pan-genomes for  $k = 19, 23, 27, 31$ . **d** Characteristics of the resulting simplitigs  
579 and unitigs for the All-dataset and Solid-dataset and  $k = 19, 23, 27, 31$ . **e** Memory footprint, index loading time and  
580 time to query 10 million  $k$ -mers using BWA (individual repetitions).

## 581 **Data**

582 All data generated or analysed during this study are included in this published article and its supplementary  
583 information files. The simplitigs of the Human genome (HG38, for  $k = 10, 11, \dots, 32$ ) and the obtained Genbank  
584 pan-genomes (for  $k = 19, 23, 27, 31$ ) are provided on Zenodo under the accessions [10.5281/zenodo.3770419](https://doi.org/10.5281/zenodo.3770419) and  
585 [10.5281/zenodo.3800713](https://doi.org/10.5281/zenodo.3800713), respectively. The code used for the analyses is provided on Github  
586 (<https://github.com/karel-brinda/simplitigs-supplementary>).

## 587 **Software**

588 ProphAsm is open source, licensed under the MIT License. The program was developed in C++ and its source code  
589 is available from Github (<http://github.com/prophyle/prophasm>). ProphAsm binaries for Linux and OS X are  
590 distributed through BioConda [83] (<https://bioconda.github.io/recipes/prophasm/README.html>). The source  
591 code of the version used in this paper was deposited in Zenodo ([10.5281/zenodo.3887035](https://doi.org/10.5281/zenodo.3887035)).

## 592 **Competing interests**

593 No competing interests.

## 594 **Funding**

595 KB and MB were partially supported by the David and Lucile Packard Foundation and NIGMS of the National  
596 Institutes of Health under award number R35GM133700. GK was partially funded by RFBR, project 20-07-00652,  
597 and joint RFBR and JSPS project 20-51-50007.

## 598 **Authors' contributions**

599 KB, MB, GK designed the study, contributed to interpretation of the results, wrote the manuscript, and approved  
600 the final manuscript. KB developed the software and performed the data analysis. KB and GK developed the theory.

## 601 **Acknowledgements**

602 The authors thank Jasmijn Baaijens and Roman Cheplyaka for careful reading and valuable comments, and Kamil  
603 Salikhov and Simone Pignotti for helpful discussions at the initial stage of this project. Portions of this research  
604 were conducted on the O2 high-performance compute cluster, supported by the Research Computing Group at  
605 Harvard Medical School.



## 606 References

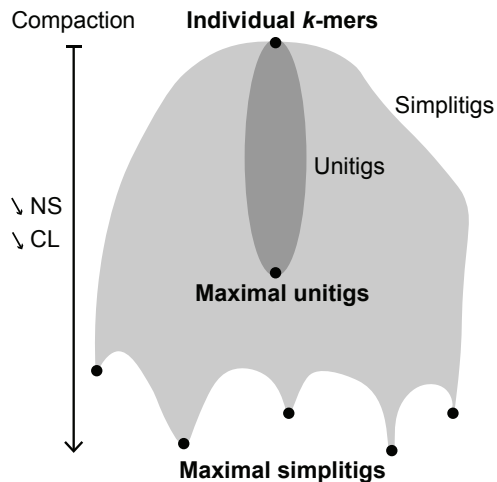
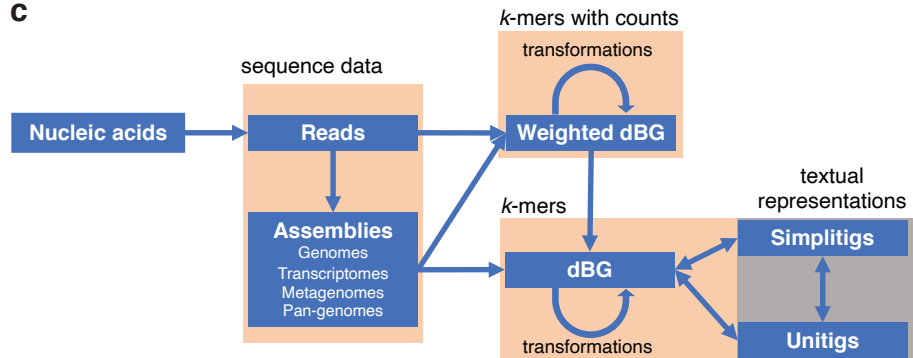
- 607 1. Needleman SB, Wunsch CD. A general method applicable to the search for similarities in the amino acid sequence  
608 of two proteins. *J Mol Biol.* 1970;48:443–53.
- 609 2. Smith TF, Waterman MS. Identification of common molecular subsequences. *J Mol Biol.* 1981;147:195–7.
- 610 3. Gotoh O. An improved algorithm for matching biological sequences. *J Mol Biol.* 1982;162:705–8.
- 611 4. Idury RM, Waterman MS. A New Algorithm for DNA Sequence Assembly. *J Comput Biol.* 1995;2:291–306.
- 612 5. Pevzner PA. 1-Tuple DNA Sequencing: Computer Analysis. *J Biomol Struct Dyn.* 1989;7:63–73.
- 613 6. Pevzner PA, Tang H, Waterman MS. An Eulerian path approach to DNA fragment assembly. *Proceedings of the*  
614 *National Academy of Sciences.* 2001;98:9748–53.
- 615 7. Zerbino DR, Birney E. Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Res.*  
616 2008;18:821–9.
- 617 8. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJM. ABySS: A parallel assembler for short read sequence  
618 data. 2009;1117–23.
- 619 9. Bankevich A, Nurk S, Antipov D, Gurevich A a., Dvorkin M, Kulikov AS, et al. SPAdes: A New Genome Assembly  
620 Algorithm and Its Applications to Single-Cell Sequencing. *J Comput Biol.* 2012;19:455–77.
- 621 10. Chikhi R, Rizk G. Space-efficient and exact de Bruijn graph representation based on a Bloom filter. *Algorithms*  
622 *Mol Biol.* 2013;8:22.
- 623 11. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node solution for large and  
624 complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics.* 2015;31:1674–6.
- 625 12. Suvorov A, Agarwala R, Lipman DJ. SKESA: strategic k-mer extension for scrupulous assemblies. *Genome*  
626 *Biol.* 2018;19:153.
- 627 13. Zielezinski A, Vinga S, Almeida J, Karlowski WM. Alignment-free sequence comparison: benefits, applications,  
628 and tools. *Genome Biol.* 2017;18:186.
- 629 14. Iqbal Z, Caccamo M, Turner I, Flicek P, McVean G. De novo assembly and genotyping of variants using colored  
630 de Bruijn graphs. *Nat Genet.* 2012;44:226–32.
- 631 15. Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, et al. Rapid antibiotic-resistance predictions from  
632 genome sequence data for *Staphylococcus aureus* and *Mycobacterium tuberculosis*. *Nat Commun. Nature*  
633 *Publishing Group;* 2015;6:10063.
- 634 16. Shajii AR, Yorukoglu D, William Yu Y, Berger B, Yu YW, Berger B. Fast genotyping of known SNPs through  
635 approximate k-mer matching. *Bioinformatics.* 2016;32:i538–44.
- 636 17. Sun C, Medvedev P. Toward fast and accurate SNP genotyping from whole genome sequencing data for bedside  
637 diagnostics. *Bioinformatics.* 2019;35:415–20.
- 638 18. Nordström KJV, Albani MC, James GV, Gutjahr C, Hartwig B, Turck F, et al. Mutation identification by direct  
639 comparison of whole-genome sequencing data from mutant and wild-type individuals using k-mers. *Nat*  
640 *Biotechnol.* 2013;31:325–30.
- 641 19. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq quantification. *Nat Biotechnol.*

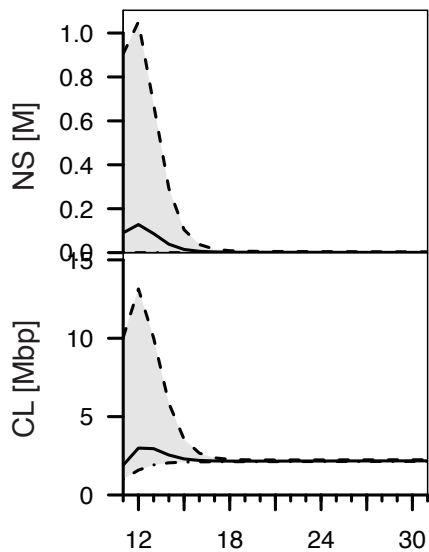
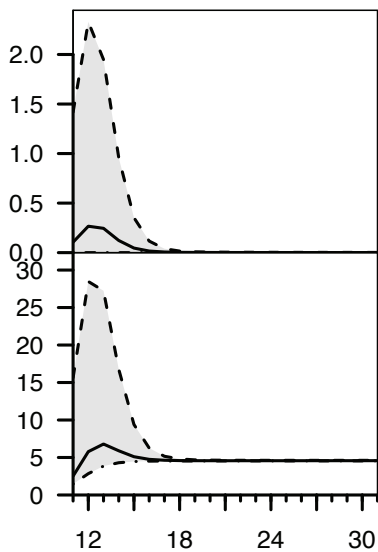
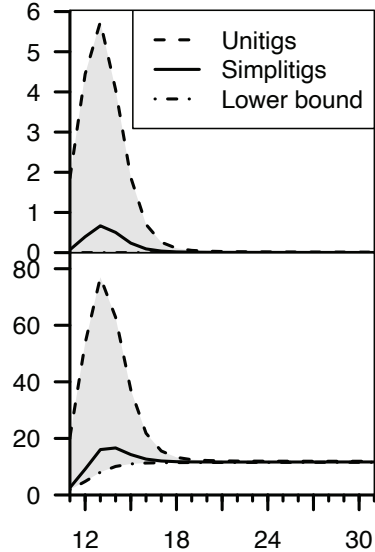
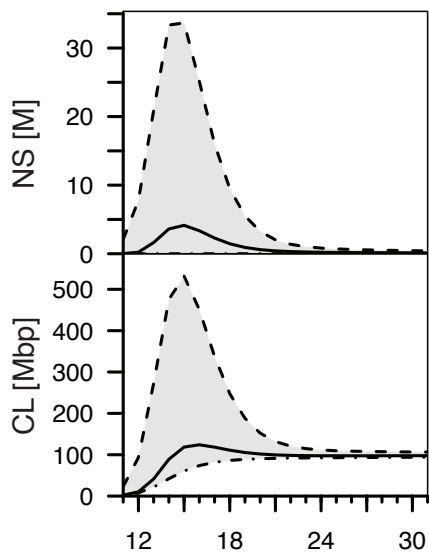
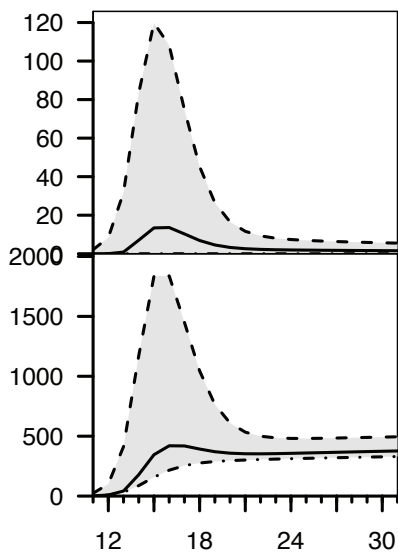
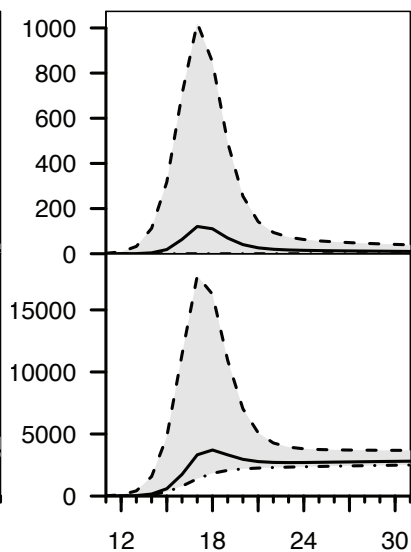
- 642 2016;34:525–7.
- 643 20. Wood DE, Salzberg SL. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome*  
644 *Biol.* 2014;15:R46.
- 645 21. Ames SK, Hysom DA, Gardner SN, Lloyd GS, Gokhale MB, Allen JE. Scalable metagenomic taxonomy  
646 classification using a reference genome database. *Bioinformatics.* 2013;29:2253–60.
- 647 22. Břinda K, Salikhov K, Pignotti S, Kucherov G. ProPhyle: An accurate, resource-frugal and deterministic DNA  
648 sequence classifier [Internet]. Zenodo; 2017. Available from: <https://zenodo.org/record/1054443>
- 649 23. Corvelo A, Clarke WE, Robine N, Zody MC. taxMaps: comprehensive and highly accurate taxonomic  
650 classification of short-read data in reasonable time. *Genome Res.* 2018;28:751–8.
- 651 24. Ye SH, Siddle KJ, Park DJ, Sabeti PC. Benchmarking Metagenomics Tools for Taxonomic Classification. *Cell.*  
652 Elsevier Inc.; 2019;178:779–94.
- 653 25. Paten B, Novak AM, Eizenga JM, Garrison E. Genome graphs and the evolution of genome inference. *Genome*  
654 *Res.* 2017;27:665–76.
- 655 26. Sirén J. Indexing Variation Graphs. 2017 Proceedings of the Ninteenth Workshop on Algorithm Engineering  
656 and Experiments (ALENEX). Philadelphia, PA: Society for Industrial and Applied Mathematics; 2017. p. 13–27.
- 657 27. Garrison E, Sirén J, Novak AM, Hickey G, Eizenga JM, Dawson ET, et al. Variation graph toolkit improves read  
658 mapping by representing genetic variation in the reference. *Nat Biotechnol.* Nature Publishing Group;  
659 2018;36:875–81.
- 660 28. Melsted P, Pritchard JK. Efficient counting of k-mers in DNA sequences using a bloom filter. *BMC*  
661 *Bioinformatics.* BioMed Central Ltd; 2011;12:333.
- 662 29. Marçais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers.  
663 *Bioinformatics.* 2011;27:764–70.
- 664 30. Deorowicz S, Debudaj-Grabysz A, Grabowski S. Disk-based k-mer counting on a PC. *BMC Bioinformatics.*  
665 2013;14:160.
- 666 31. Rizk G, Lavenier D, Chikhi R. DSK: k-mer counting with very low memory usage. *Bioinformatics.*  
667 2013;29:652–3.
- 668 32. Roy RS, Bhattacharya D, Schliep A. Turtle: Identifying frequent k-mers with cache-efficient algorithms.  
669 *Bioinformatics.* 2014;30:1950–7.
- 670 33. Crusoe MR, Alameldin HF, Awad S, Boucher E, Caldwell A, Cartwright R, et al. The khmer software package:  
671 enabling efficient nucleotide sequence analysis. *F1000Res.* 2015;1–12.
- 672 34. Deorowicz S, Kokot M, Grabowski S, Debudaj-Grabysz A. KMC 2: fast and resource-frugal k-mer counting.  
673 *Bioinformatics.* 2015;31:1569–76.
- 674 35. Kokot M, Dlugosz M, Deorowicz S. KMC 3: counting and manipulating k-mer statistics. *Bioinformatics.*  
675 2017;33:2759–61.
- 676 36. Pandey P, Bender MA, Johnson R, Patro R, Berger B. Squeakr: an exact and approximate k-mer counting  
677 system. *Bioinformatics.* 2018;34:568–75.
- 678 37. Chikhi R, Limasset A, Jackman S, Simpson JT, Medvedev P. On the Representation of De Bruijn Graphs. *J*  
679 *Comput Biol.* 2015;22:336–52.
- 680 38. Chikhi R, Limasset A, Medvedev P. Compacting de Bruijn graphs from sequencing data quickly and in low

- 681 memory. *Bioinformatics*. 2016;32:i201–8.
- 682 39. Břinda K, Callendrello A, Ma KC, MacFadden DR, Charalampous T, Lee RS, et al. Rapid inference of antibiotic  
683 resistance and susceptibility by genomic neighbour typing. *Nature Microbiology* [Internet]. 2020; Available from:  
684 <http://dx.doi.org/10.1038/s41564-019-0656-6>
- 685 40. Marschall T, Marz M, Abeel T, Dijkstra L, Dutilh BE, Ghaffaari A, et al. Computational pan-genomics: status,  
686 promises and challenges. *Brief Bioinform*. 2016;184:bbw089.
- 687 41. Grad YH, Harris SR, Kirkcaldy RD, Green AG, Marks DS, Bentley SD, et al. Genomic Epidemiology of  
688 Gonococcal Resistance to Extended-Spectrum Cephalosporins, Macrolides, and Fluoroquinolones in the United  
689 States, 2000–2013. *J Infect Dis*. 2016;214:1579–87.
- 690 42. Croucher NJ, Finkelstein JA, Pelton SI, Mitchell PK, Lee GM, Parkhill J, et al. Population genomics of  
691 post-vaccine changes in pneumococcal epidemiology. *Nat Genet*. Nature Publishing Group; 2013;45:656–63.
- 692 43. Croucher NJ, Finkelstein JA, Pelton SI, Parkhill J, Bentley SD, Lipsitch M, et al. Population genomic datasets  
693 describing the post-vaccine evolutionary epidemiology of *Streptococcus pneumoniae*. *Scientific data*.  
694 2015;2:150058.
- 695 44. Mäkinen V, Belazzougui D, Cunial F, Tomescu AI. *Genome-Scale Algorithm Design*. Cambridge University  
696 Press; 2015.
- 697 45. Ferragina P, Manzini G. Opportunistic data structures with applications. *Proceedings 41st Annual Symposium*  
698 *on Foundations of Computer Science*. IEEE Comput. Soc; 2000. p. 390–8.
- 699 46. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*.  
700 *Narnia*; 2009;25:1754–60.
- 701 47. Li H. Exploring single-sample SNP and INDEL calling with whole-genome de novo assembly. *Bioinformatics*.  
702 2012;28:1838–44.
- 703 48. Merchant S, Wood DE, Salzberg SL. Unexpected cross-species contamination in genome sequencing projects.  
704 *PeerJ*. 2014;2:e675.
- 705 49. Lu J, Salzberg SL. Removing contaminants from databases of draft genomes. *PLoS Comput Biol*.  
706 2018;14:e1006277.
- 707 50. Steinegger M, Salzberg SL. Terminating contamination: large-scale search identifies more than 2,000,000  
708 contaminated entries in GenBank [Internet]. *bioRxiv*. 2020. p. 2020.01.26.920173. Available from:  
709 <https://www.biorxiv.org/content/10.1101/2020.01.26.920173v1>
- 710 51. Goig GA, Blanco S, Garcia-Basteiro AL, Comas I. Contaminant DNA in bacterial sequencing experiments is a  
711 major source of false genetic variability. *BMC Biol*. 2020;18:24.
- 712 52. Lin Y, Yuan J, Kolmogorov M, Shen MW, Chaisson M, Pevzner PA. Assembly of long error-prone reads using de  
713 Bruijn graphs. *Proc Natl Acad Sci U S A*. 2016;113:E8396–405.
- 714 53. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs. *Nat*  
715 *Biotechnol*. 2019;37:540–6.
- 716 54. Břinda K. Novel computational techniques for mapping and classification of Next-Generation Sequencing data  
717 [Internet]. 2016. Available from: <https://hal.archives-ouvertes.fr/tel-01484198/>
- 718 55. Li H. Revisiting hash table performance [Internet]. *Attractive Chaos*. 2018 [cited 2020 May 1]. Available from:  
719 <https://attractivechaos.wordpress.com/2018/01/13/revisiting-hash-table-performance/>
- 720 56. Pan T, Nihalani R, Aluru S. Fast de Bruijn Graph Compaction in Distributed Memory Environments.

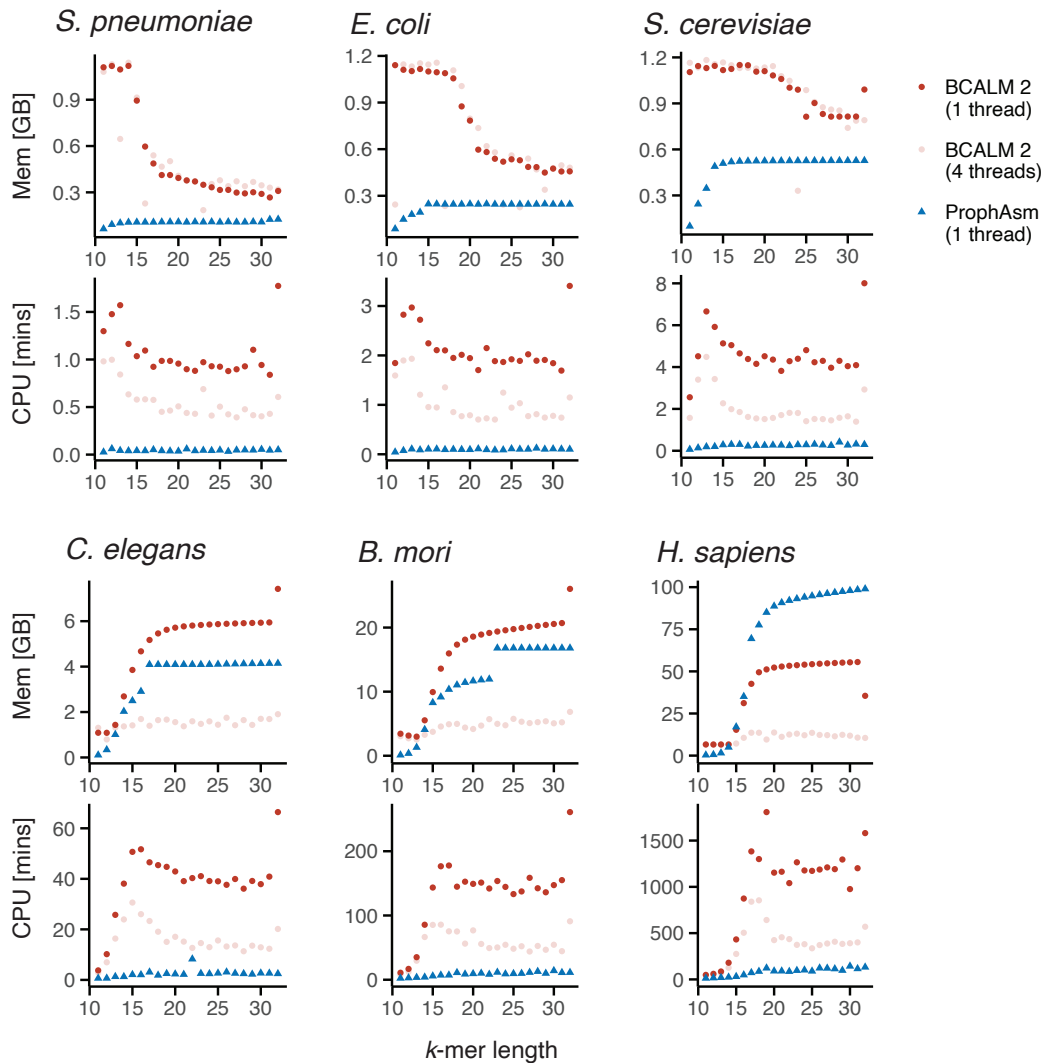
- 721 IEEE/ACM Trans Comput Biol Bioinform. 2018;1–1.
- 722 57. Guo H, Fu Y, Gao Y, Li J, Wang Y, Liu B. deGSM: memory scalable construction of large scale de Bruijn Graph.  
723 IEEE/ACM Trans Comput Biol Bioinform. 2019;1–1.
- 724 58. Nasko DJ, Koren S, Phillippy AM, Treangen TJ. RefSeq database growth influences the accuracy of k-mer-based  
725 lowest common ancestor species identification. *Genome Biol.* 2018;19:165.
- 726 59. Bowe A, Onodera T, Sadakane K, Shibuya T. Succinct de Bruijn Graphs. In: Raphael B, Tang J, editors.  
727 Algorithms in Bioinformatics. Berlin, Heidelberg: Springer Berlin Heidelberg; 2012. p. 225–35.
- 728 60. Holley G, Wittler R, Stoye J. Bloom Filter Trie: an alignment-free and reference-free data structure for  
729 pan-genome storage. *Algorithms Mol Biol. BioMed Central*; 2016;11:3.
- 730 61. Solomon B, Kingsford C. Fast search of thousands of short-read sequencing experiments. *Nat Biotechnol.*  
731 2016;34:300–2.
- 732 62. Muggli MD, Bowe A, Noyes NR, Morley PS, Belk KE, Raymond R, et al. Succinct colored de Bruijn graphs.  
733 *Bioinformatics.* 2017;33:3181–7.
- 734 63. Sun C, Harris RS, Chikhi R, Medvedev P. AllSome Sequence Bloom Trees. *J Comput Biol.* 2018;25:467–79.
- 735 64. Pandey P, Almodaresi F, Bender MA, Ferdman M, Johnson R, Patro R. Mantis: A Fast, Small, and Exact  
736 Large-Scale Sequence-Search Index. *Cell Syst.* 2018;7:201–7.e4.
- 737 65. Yu Y, Liu J, Liu X, Zhang Y, Magner E, Lehnert E, et al. SeqOthello: querying RNA-seq experiments at scale.  
738 *Genome Biol.* 2018;19:167.
- 739 66. Almodaresi F, Sarkar H, Srivastava A, Patro R. A space and time-efficient index for the compacted colored de  
740 Bruijn graph [Internet]. *Bioinformatics.* 2018. p. i169–77. Available from:  
741 <http://dx.doi.org/10.1093/bioinformatics/bty292>
- 742 67. Harris RS, Medvedev P. Improved representation of sequence Bloom trees [Internet]. *Bioinformatics.* 2019.  
743 Available from: <http://dx.doi.org/10.1093/bioinformatics/btz662>
- 744 68. Holley G, Melsted P. Bifrost – Highly parallel construction and indexing of colored and compacted de Bruijn  
745 graphs [Internet]. *bioRxiv.* 2019. p. 695338. Available from: <https://www.biorxiv.org/content/10.1101/695338v2>
- 746 69. Bradley P, den Bakker HC, Rocha EPC, McVean G, Iqbal Z. Ultrafast search of all deposited bacterial and viral  
747 genomic data. *Nat Biotechnol. Springer US*; 2019;37:152–9.
- 748 70. Bingmann T, Bradley P, Gauger F, Iqbal Z. COBS: A Compact Bit-Sliced Signature Index. *String Processing and*  
749 *Information Retrieval. Springer International Publishing*; 2019. p. 285–303.
- 750 71. Zielezinski A, Girgis HZ, Bernard G, Leimeister C-A, Tang K, Dencker T, et al. Benchmarking of alignment-free  
751 sequence comparison methods. *Genome Biol.* 2019;20:144.
- 752 72. Pandey P, Bender MA, Johnson R, Patro R. deBGR: an efficient and near-exact representation of the weighted  
753 de Bruijn graph. *Bioinformatics.* 2017;33:i133–41.
- 754 73. Rahman A, Medvedev P. Representation of k-mer sets using spectrum-preserving string sets [Internet]. *bioRxiv.*  
755 2020 [cited 2020 Jan 20]. p. 2020.01.07.896928. Available from:  
756 <https://www.biorxiv.org/content/10.1101/2020.01.07.896928v1.abstract>
- 757 74. Rowe WPM. When the levee breaks: a practical guide to sketching algorithms for processing the flood of  
758 genomic data. *Genome Biol. Genome Biology*; 2019;20:199.
- 759 75. Elworth RAL, Wang Q, Kota PK, Barberan CJ, Coleman B, Balaji A, et al. To Petabytes and beyond: recent  
760 advances in probabilistic and signal processing algorithms and their application to metagenomics. *Nucleic Acids*

- 761 Res [Internet]. 2020; Available from: <http://dx.doi.org/10.1093/nar/gkaa265>
- 762 76. Pinho AJ, Ferreira PJSG, Garcia SP, Rodrigues JMOS. On finding minimal absent words. *BMC Bioinformatics*.  
763 2009;10:137.
- 764 77. Gallant J, Maier D, Astorer J. On finding minimal length superstrings. *J Comput System Sci*. 1980;20:50–8.
- 765 78. Manuel P. Revisiting path-type covering and partitioning problems [Internet]. arXiv [math.CO]. 2018. Available  
766 from: <http://arxiv.org/abs/1807.10613>
- 767 79. Medvedev P, Brudno M. Maximum likelihood genome assembly. *J Comput Biol*. 2009;16:1101–16.
- 768 80. Köster J, Rahmann S. Snakemake-a scalable bioinformatics workflow engine. *Bioinformatics*. 2012;28:2520–2.
- 769 81. Grad Y. Data for “Genomic Epidemiology of Gonococcal Resistance to Extended-Spectrum Cephalosporins,  
770 Macrolides, and Fluoroquinolones in the United States, 2000-2013” [Internet]. Zenodo; 2019. Available from:  
771 <https://zenodo.org/record/2618836>
- 772 82. Homer N. DWGSIM: Whole Genome Simulator for Next-Generation Sequencing. GitHub repository. 2010;
- 773 83. Grüning B, Dale R, Sjödin A, Chapman BA, Rowe J, Tomkins-Tinch CH, et al. Bioconda: sustainable and  
774 comprehensive software distribution for the life sciences. *Nat Methods*. 2018;15:475–6.
- 775

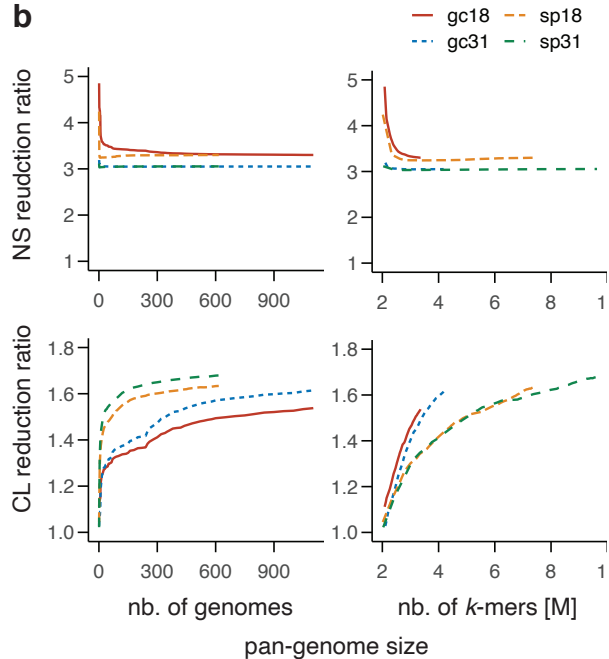
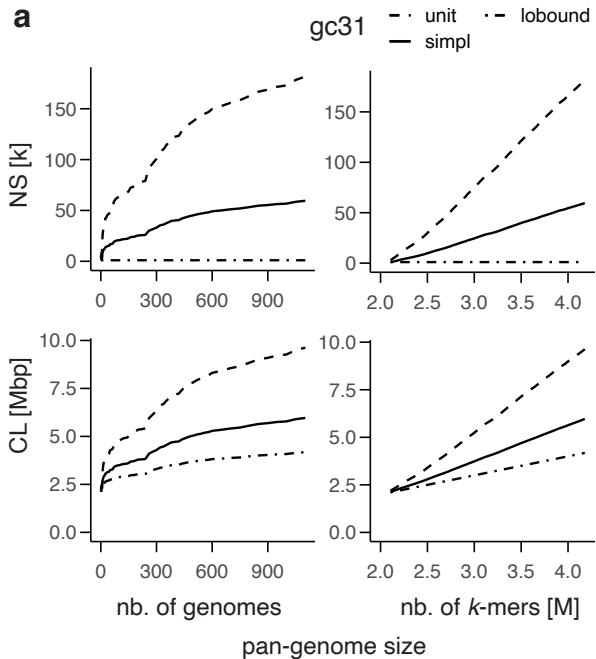
**a****Individual *k*-mers:** ACG, CGA, GAA, AAG, AGC, CGT, GTA, TAG, CGG**Maximal unitigs:** ACG, CGAAG, CGTAG, AGC, CGG**Maximal simplitigs:** ACGAAGC, CGTAG, CGG**b****c**

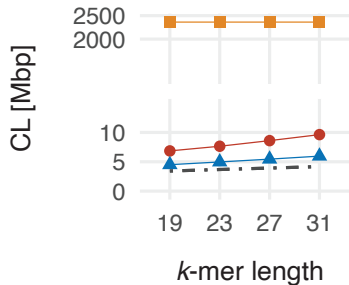
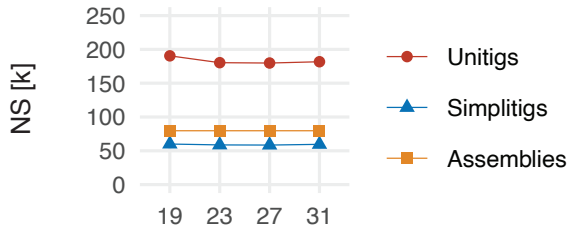
**a** *S. pneumoniae***b** *E. coli***c** *S. cerevisiae***d** *C. elegans***e** *B. mori***f** *H. sapiens*

*k*-mer length







**a****b**